

DEEP BOLTZMANN MACHINES IN ESTIMATION OF DISTRIBUTION ALGORITHMS FOR COMBINATORIAL OPTIMIZATION

MALTE PROBST AND FRANZ ROTHLAUF

ABSTRACT. Estimation of Distribution Algorithms (EDAs) require flexible probability models that can be efficiently learned and sampled. Deep Boltzmann Machines (DBMs) are generative neural networks with these desired properties. We integrate a DBM into an EDA and evaluate the performance of this system in solving combinatorial optimization problems with a single objective. We compare the results to the Bayesian Optimization Algorithm. The performance of DBM-EDA was superior to BOA for difficult additively decomposable functions which are separable, i.e., concatenated deceptive traps of higher order. For most other benchmark problems, DBM-EDA cannot clearly outperform BOA, or other neural network-based EDAs. In particular, it often yields optimal solutions for a subset of the runs (with fewer evaluations than BOA), but is unable to provide reliable convergence to the global optimum competitively. At the same time, the model building process is fast, but computationally more expensive than that of other EDAs using probabilistic models from the neural network family, such as DAE-EDA.

1. INTRODUCTION

Estimation of Distribution Algorithms (EDAs) [12, 8] are metaheuristics for combinatorial and continuous non-linear optimization. They maintain a population of candidate solutions to the optimization problem at hand (see Algorithm 1). First, they select solutions with a high quality from the population. Subsequently, they build a model that approximates the probability distribution of these solutions. Then, new candidate solutions are sampled from the model. The EDA then starts over by selecting the next set of good solutions from the new candidate solutions and the previous selection.

In order to be suitable for an EDA, a model therefore has to fulfill certain criteria:

- It must be able to approximate the probability distribution of the selected individuals.
- It must be able to sample new solutions from this probability distribution, serving as candidate solutions for the next EDA generation.
- Both learning and sampling should be efficient. That is, the computational time required to train and sample the model should be tractable both the number of variables, and in the number of training examples.

Previous work has shown that generative neural networks can lead to competitive performance. [22] use a Restricted Boltzmann Machine (RBM) in an EDA and show that RBM-EDA can achieve competitive performance to state-of-the-art EDAs,

Key words and phrases. Combinatorial Optimization, Heuristics, Evolutionary Computation, Estimation of Distribution Algorithm, Neural Networks, Deep Boltzmann Machine .

especially in terms of computational complexity of the CPU time. [20] use another type of neural network, a Denoising Autoencoder (DAE) in an EDA. DAE-EDA achieves superior performance when used on problems which can be decomposed into independent subproblems. In general, neural-network inspired probabilistic models can often be parallelized on massively parallel systems such as graphics processing units (GPUs) [21].

In this paper, we focus on Deep Boltzmann Machines (DBMs) [23]. DBMs are *deep* models in the sense that they use multiple layers \mathbf{h}^j , $j = 1 \dots d$ of hidden (latent) neurons¹. A DBM models the joint probability distribution $P(\mathbf{v}, \mathbf{h}^1, \dots, \mathbf{h}^d)$ of the training data \mathbf{v} and the hidden neurons.

Deep architectures are particularly interesting, because they are able to model problems on multiple layers of abstraction. A deep model is usually composed of multiple layers of computational units (e.g., neurons). The concepts modeled by each layer becomes more abstract with the layer’s depth [2, 9]. An intuitive example is a deep neural network that learns to model images of faces [10]: Neurons on the first hidden layer learn to model individual edges and other shapes. Units on deeper layers compose these edges to form higher-level features, like noses or eyes. Again, by combining these mid-level representations, neurons in the deepest layers can compose complete faces. Many real-world problems like image classification possess this kind of hierarchic structure with various layers of abstraction. Deep models have recently gained much attention, as they were able to yield superior results for various real-world problem domains [9].

We implement a DBM and use it within an EDA to solve combinatorial optimization problems. We test DBM-EDA on the simple onemax problem, concatenated deceptive trap functions, NK landscapes and the HIFF function. We compare the results the state-of-the-art multivariate Bayesian Optimization Algorithm (BOA, see [17, 15]), and publish the source code of all experiments²

Section 2, introduces DBMs. Section 3 describes benchmark problems, experimental setup, and presents the results. We discuss the results and conclude the paper in section 4.

2. DEEP BOLTZMANN MACHINES

DBMs are special types of Boltzmann Machines, one of the fundamental types of neural networks [1, 23]. A DBM has a visible layer $\mathbf{v} \in [0, 1]^n$ and d hidden

¹We use the following notation: x denotes a scalar value, \mathbf{x} denotes a vector of scalars, \mathbf{X} denotes a matrix of scalars

²See <https://github.com/wohnjayne/eda-suite/> for the complete source code

Algorithm 1 Estimation of Distribution Algorithm

- 1: **Initialize** Population P
 - 2: **while** not converged **do**
 - 3: $P_{parents} \leftarrow$ **Select** high-quality solutions from P based on their fitness
 - 4: $M \leftarrow$ **Build** a model estimating the (joint) probability distribution of $P_{parents}$
 - 5: $P_{candidates} \leftarrow$ **Sample** new candidate solutions from M
 - 6: $P \leftarrow P_{parents} \cup P_{candidates}$
 - 7: **end while**
-

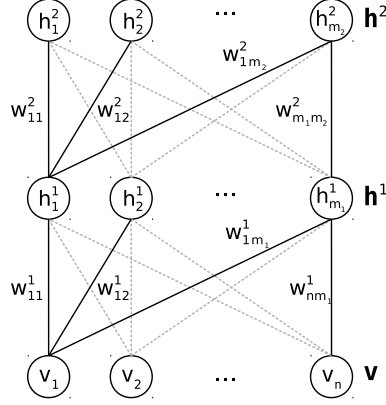


FIGURE 1. A Deep Boltzmann Machine with two hidden layers $\mathbf{h}^1, \mathbf{h}^2$ as a graph. The visible neurons v_i ($i \in 1..n$) can hold a data vector of length n from the training data. In the EDA context, \mathbf{v} represents decision variables. The hidden neurons h_j^1 ($j \in 1 \dots m_1$) and h_k^2 ($k \in 1 \dots m_2$) represent m_1 first-level and m_2 second-level features, respectively.

layers $\mathbf{h}^j \in [0, 1]^{m_j}$. Each layer consists of binary units, and learns a non-linear representation of the data on the layer below. Hence, upper layers will learn more abstract concepts.

Here, we focus on DBMs with two hidden layers \mathbf{h}^1 and \mathbf{h}^2 (see Figure 1). The layers of neurons are connected by symmetric weights. The weight matrix \mathbf{W}^1 of size $n * m_1$ connects the visible layer to the first hidden layer. Weight w_{ij}^1 therefore connects v_i to h_j^1 . Accordingly, weight matrix \mathbf{W}^2 of size $m_1 * m_2$ connects hidden layer \mathbf{h}^1 to hidden layer \mathbf{h}^2 . There are no connections within the layers.

From Boltzmann Machines, DBMs inherit the concept of a scalar *energy* associated with each configuration of its neurons. The energy of the state $\{\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2\}$ is defined as

$$(1) \quad E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \theta) = -\mathbf{v}^\top \mathbf{W}^1 \mathbf{h}^1 - \mathbf{h}^1^\top \mathbf{W}^2 \mathbf{h}^2$$

where $\theta = \{\mathbf{W}^1, \mathbf{W}^2\}$ are the model's parameters³. The probability of a particular configuration of the visible neurons \mathbf{v} under the model is

$$(2) \quad P(\mathbf{v}; \theta) = \frac{1}{Z(\theta)} \sum_{\mathbf{h}^1, \mathbf{h}^2} \exp(-E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \theta)).$$

$Z(\theta)$ is the partition function which normalizes the probability by summing over all possible configurations.

³We omit the bias terms for brevity, see [24].

The *conditional* probabilities of the neurons, given the activations of their neighboring layers, is easier to calculate. Each neuron v_i calculates its conditional probability to be active as

$$(3) \quad P(v_i = 1 | \mathbf{h}^1) = \text{sigm}(\sum_j h_j^1 * w_{ij}^{v h_1}),$$

with $\text{sigm}(x)$ being the logistic function $\text{sigm}(x) = \frac{1}{1+e^{-x}}$. The neurons in hidden layer h^1 calculate their probability to be active as

$$(4) \quad P(h_j^1 = 1 | \mathbf{v}, \mathbf{h}^2) = \text{sigm}(\sum_i v_i * w_{ij}^1 + \sum_k h_k^2 * w_{jk}^2),$$

and the neurons in hidden layer h^2 calculate their probability to be active as

$$(5) \quad P(h_k^2 = 1 | \mathbf{h}^1) = \text{sigm}(\sum_j h_j^1 * w_{jk}^2).$$

Algorithm 2 Pseudo code for pre-training a DBM with RBMs

- 1: **Initialize** $\mathbf{W}^1, \mathbf{W}^2$ to small, random values
 - 2: $\mathbf{W}^1 \leftarrow$ Train RBM₁ modeling $P(\mathbf{v}, \mathbf{h}^1; \mathbf{W}^1)$,
 - 3: using P_{Parents} as training set
 - 4: $\mathbf{W}^2 \leftarrow$ Train RBM₂ modeling $P(\mathbf{h}^1, \mathbf{h}^2; \mathbf{W}^2)$,
 - 5: using samples from $P(\mathbf{H}^1 | P_{\text{Parents}}; \mathbf{W}^1)$ as defined by RBM₁ as training set
-

2.1. Training a DBM. In the training phase, the parameters θ of the DBM have to be adjusted such that the model approximates the probability distribution of the training data. This could, in principle, be done by using the general training procedure of Boltzmann Machines [1]. However, this would be computationally intractable. Hence, a greedy, layer-wise pretraining is used to initialize the parameters to some sensible values, and subsequently perform parameter fine-tuning on the complete DBM.

During the pretraining phase, we consider the DBM to be a stack of Restricted Boltzmann Machines (RBMs) (see Algorithm 2). RBMs are similar to DBMs, but only have a single layer of hidden neurons. Contrastive divergence is a tractable learning algorithm to train RBMs [4]. Specifically, we consider \mathbf{v} and \mathbf{h}^1 to be an RBM, with \mathbf{W}^1 as its parameters. We then train the resulting RBM to model the probability distribution of the training data, using contrastive divergence (see e.g. [4, 5] for details). Subsequently, we train a second RBM, consisting of $\mathbf{h}^1, \mathbf{h}^2$ and \mathbf{W}^2 . The training data for the second RBM consists of the activations of the first RBM's hidden layer \mathbf{h}^1 on the original training data⁴.

Once both RBMs have been trained, we use their parameters \mathbf{W}^1 and \mathbf{W}^2 as an initialization for the parameters θ of a single DBM. Note that the pretraining does *not* adjust the weights s.t. the DBM's conditional probability distribution $P(\mathbf{v} | \mathbf{h}^1, \mathbf{h}^2; \theta)$ approximates the probability distribution of the training data.

In order to achieve this, we fine-tune the DBM (see Algorithm 3). This can be done using gradient descent algorithm to modify θ . The general idea of the

⁴There is a small adjustment in the training procedure for the RBMs. In practice, the weight matrices \mathbf{W}^1 and \mathbf{W}^2 are multiplied by constant factors, see [24] for details.

Algorithm 3 Pseudo code for fine-tuning a DBM

```

1: Set  $0 < \alpha < 1$ , e.g.  $\alpha = 0.1$ 
2: Initialize DBM's parameters with pre-trained  $\mathbf{W}^1$  and  $\mathbf{W}^2$ 
3: Initialize fantasy particles  $\hat{\mathbf{v}}$ ,  $\hat{\mathbf{h}}^1$  and  $\hat{\mathbf{h}}^2$  randomly
4: while not converged do
5:   for each example in the training set do
6:     — Positive phase —
7:      $\mathbf{v} \leftarrow$  set  $\mathbf{v}$  to the current training example
8:      $\mathbf{h}^1 \leftarrow$  set to  $P(\mathbf{h}^1|\mathbf{v})$ ,
9:       ignoring input from  $\mathbf{h}^2$ 
10:     $\mathbf{h}^2 \leftarrow$  set to  $P(\mathbf{h}^2|\mathbf{h}^1)$ 
11:    run mean field approximation:
12:    for a small number of iterations (e.g. ten) do
13:       $\mathbf{h}^1 \leftarrow$  set to  $P(\mathbf{h}^1|\mathbf{v}, \mathbf{h}^2)$ 
14:       $\mathbf{h}^2 \leftarrow$  set to  $P(\mathbf{h}^2|\mathbf{h}^1)$ 
15:    end for
16:    — Negative phase —
17:    for a small number of iterations (e.g. five) do
18:       $\hat{\mathbf{h}}^1 \leftarrow$  activate stochastically with  $P(\hat{\mathbf{h}}^1|\hat{\mathbf{v}}, \hat{\mathbf{h}}^2)$ 
19:       $\hat{\mathbf{h}}^2 \leftarrow$  activate stochastically with  $P(\hat{\mathbf{h}}^2|\hat{\mathbf{h}}^1)$ 
20:       $\hat{\mathbf{v}} \leftarrow$  activate stochastically with  $P(\hat{\mathbf{v}}|\hat{\mathbf{h}}^1)$ 
21:    end for
22:    — Calculate and apply gradient —
23:     $\delta_{\text{pos}} = \frac{\partial E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \theta)}{\partial \theta}$ 
24:     $\delta_{\text{neg}} = -\frac{\partial E(\hat{\mathbf{v}}, \hat{\mathbf{h}}^1, \hat{\mathbf{h}}^2; \theta)}{\partial \theta}$ 
25:     $\theta := \theta - \alpha * (\delta_{\text{pos}} + \delta_{\text{neg}})$ 
26:  end for
27: end while

```

algorithm is as follows: In order to increase the probability of a training data point \mathbf{v} under the model, the energy $E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \theta)$ of this configuration has to be decreased (see Equation 2). At the same time, the probability of all other configurations has to be lowered, by increasing their energies contained in the partition function Z .

The gradient hence contains two terms. The first term (positive gradient) increases the probability of the current sample under the model, the second term (negative gradient) decreases the probabilities of all other configurations in Z .

The negative gradient is calculated by running a separate Markov chain of *fantasy particles* in the model. Specifically, all neurons are first initialized randomly. Then, the following two steps are repeated: First, \mathbf{h}^1 is activated stochastically, with probabilities $P(\mathbf{h}^1|\mathbf{v}, \mathbf{h}^2)$ as in Equation 4. Then, \mathbf{v} and \mathbf{h}^2 are activated stochastically, with $P(\mathbf{v}|\mathbf{h}^1)$ as in Equation 3, and $P(\mathbf{h}^2|\mathbf{h}^1)$ as in Equation 5. If the parameter updates are small enough, and the Markov chain is allowed to run a couple of steps between each update (e.g. five steps), the samples will come from the chain's equilibrium distribution. The current state of the fantasy particle is then used to calculate the negative gradient.

The positive term could be approximated by running the same Markov chain as above, but with \mathbf{v} clamped to the current training example. However, this would

result in running the Markov chain for many steps, for each training example. Instead, the positive gradient can be approximated using a mean-field approach. First, neurons in \mathbf{v} are set according to the current training example. Then, $P(\mathbf{h}^1)$ and, subsequently, $P(\mathbf{h}^2|P(\mathbf{h}^1))$ are calculated, as in Equations 4 and 5. Then, for a small number of steps (e.g. ten steps), $P(\mathbf{h}^1|\mathbf{v}, P(\mathbf{h}^2))$ and $P(\mathbf{h}^2|P(\mathbf{h}^1))$ are calculated repeatedly. The resulting configuration of \mathbf{v} , $P(\mathbf{h}^1)$ and $P(\mathbf{h}^2)$ is treated as a positive example, and used for the calculation of the positive gradient. The DBM's parameters θ are then updated in the direction of the total gradient.

For a more detailed description of the training algorithm, see [24].

2.2. Sampling a DBM. The DBM can be sampled by initializing all neurons to random values, and running the same sampling chain used to retrieve fantasy particles for the negative gradient. In the case of DBM-EDA, we initialize v with the P_{Parents} , and run the chain for 25 iterations.

3. EXPERIMENTS

We evaluate DBM-EDA using a set of standard benchmark problems. We compare the results of DBM-EDA to those of the state-of-the-art multivariate BOA.

3.1. Test Problems. We evaluate DBM-EDA on onemax, concatenated deceptive traps, NK landscapes and the HIFF function. All four are standard benchmark problems. Their difficulty depends on the problem size, i.e., problems with more decision variables are more difficult. Furthermore, the difficulty of concatenated deceptive trap functions and NK landscapes is tunable by a parameter. Apart from the simple onemax problem, all problems are composed of subproblems, which are either deceptive (traps), overlapping (NK landscapes), or hierarchical (HIFF), and therefore multimodal.

The onemax problem assigns a binary solution \mathbf{x} of length l a fitness value $f = \sum_{i=1}^l x_i$, i.e., the fitness of \mathbf{x} is equal to the number of ones in \mathbf{x} . The onemax function is rather simple. It is unimodal and can be solved by a deterministic hill climber.

Concatenated deceptive traps are tunably hard, yet separable test problems [3]. Here, a solution vector \mathbf{x} is divided into l subsets of size k , with each one being a deceptive trap. Within a trap, all bits are dependent on each other but independent of all other bits in \mathbf{x} . Thus, the fitness contribution of the traps can be evaluated separately and the total fitness of the solution vector is the sum of these terms. In particular, the assignment $\mathbf{a} = \mathbf{x}_{i:i+k-1}$ (i.e., the k bits from x_i to x_{i+k-1})⁵ leads to a fitness contribution F_l as

$$F_l(\mathbf{a}) = \begin{cases} k & \text{if } \sum_i a_i = k, \\ k - (\sum_i a_i + 1) & \text{otherwise.} \end{cases}$$

In other words, the fitness of a single trap increases with the number of zeros, except for the optimum of all ones.

NK landscapes are defined by two parameters n and k and n fitness components $f_i, i \in \{1 \dots, n\}$ [6]. A solution vector \mathbf{x} consists of n bits. The bits are assigned to n overlapping subsets, each of size $k + 1$. The fitness of a solution is the sum of n

⁵The k variables assigned to trap l do not have to be adjacent, but can be at any position in x .

fitness components. Each component f_i depends on the value of the corresponding variable x_i as well as k other variables. Each f_i maps each possible configurations of its $k + 1$ variables to a fitness value. The overall fitness function is

$$f(\mathbf{x}) = 1/n \sum_{i=1}^n f_i(x_i, x_{i1}, \dots, x_{iK}).$$

Each decision variable usually influences several f_i . These dependencies between subsets make NK landscapes non-separable, i.e., in general, we cannot solve the subproblems independently. The problem difficulty increases with k . $k = 0$ is a special case where all decision variables are independent and the problem reduces to a unimodal onemax. We use instances of NK landscapes with known optima from [16].

The Hierarchical If-and-only-if (HIFF) function [27] is defined for solutions vectors of length $n = 2^l$ where $l \in \mathbb{N}$ is the number of layers of the hierarchy. It uses a mapping function M and a contribution function C , both of which take two inputs. The mapping function takes each of the $n/2$ blocks of two neighboring variables of level $l = 1$, and maps them onto a single symbol each. An assignment of 00 is mapped to 0, 11 is mapped to 1 and everything else is mapped to the null symbol '-'. The concatenation of M 's outputs on level l is used as M 's input for the next level $l + 1$ of the hierarchy, i.e., if level $l = 1$ has n variables, level $l = 2$ has $n/2$ variables. On each level, C assigns a fitness to each block of two variables. The assignments 00 and 11 are both mapped to 2^l , everything else to 0. The total fitness is the sum of all blocks' contributions on all levels. In other words, a block contributes to the fitness on the current level if both variables in a block have the same assignment. However, only if neighboring blocks agree on the assignment, they will contribute to the fitness on the next level, which is why HIFF is a difficult problem. HIFF has two global optima, the string of all ones, and the string of all zeros.

3.2. Experimental Setup. We use several instances of the test problems. For each instance and algorithm, we test multiple population sizes between 100 and 16,000⁶.

We run 20 instances for each population size. In each run, the EDAs are allowed to run for up to 150 generations. We terminate a run if there is no improvement in the best solution for more than 50 generations. These settings make it very unlikely that a run is terminated prematurely, i.e., before convergence. Both DBM-EDA and BOA use tournament selection without replacement of size two [11]. Note that all test problems, with the exception of NK landscapes, have the string of all ones as their global optimum, for any problem size. To avoid any possible model-induced bias towards solutions with ones or zeros, we generate a random matrix $R \in [0, 1]^{n \times m}$ of ones and zeros for each run. In each generation, we apply the following operations. Before training, we set $\text{trainingData} = P_{\text{parents}} \oplus R$, with \oplus being a logical XOR. After sampling we set $P_{\text{candidates}} = \text{modelSamples} \oplus R$. and after sampling. These operations are transparent to correlations between variables.

We use standard values for all hyper-parameters governing the DBM's learning and sampling procedures. All hyper-parameters, and further details of the learning process such as momentum or weight decay are available in a configuration file along with the source code (see git repository on github.com).

⁶popsize $\in \{100; 200; 300; 400; 500; 1,000; 1,500; \text{ and } 2,000 \text{ to } 16,000 \text{ (increment } 1000)\}$

The algorithms are implemented in Matlab/Octave and executed using Octave V3.2.4 on a single core of an AMD Opteron 6272 processor with 2,100 MHz. For the DBM, we used the source code provided by [23].

3.3. Results. Table 1 shows results for DBM-EDA and BOA on the Onemax problem (50, 75, 100, and 150 bit problem size), concatenated 4-Traps (40, 60, and 80 bit), concatenated 5-Traps (25, 50, and 100 bit), NK landscapes with $k \in \{4, 5\}$ (30 and 34 bit, two instances each) and the HIFF function (64 and 128 bit).

We report the population sizes, the average number of unique fitness evaluations, and the average CPU times that each algorithm needed to solve the respective problem instance to optimality in at least 50% and 90% of the runs.

First, we concentrate on the number of fitness evaluations, and on the results for solving at least 50% of the runs (left three result columns of table 1). For the simple onemax problem, DBM-EDA needs slightly less fitness evaluations than BOA. For the concatenated deceptive traps, the results are mixed. BOA finds the optimal solutions to the 4-Trip problems faster than DBM-EDA, while DBM-EDA seems to be more competitive for the harder 5-Traps problem. For the NK landscapes, DBM-EDA needs less fitness evaluations than BOA in six out of eight instances. On the HIFF problem, DBM-EDA’s performance is clearly inferior to BOA: it needs much more fitness evaluations on the 64 bit instance, and is unable to find the optimal solution in at least 50% of runs for the 128 bit instance. This is surprising, given that HIFF is a hierarchical problem, and the DBM is a hierarchical model. In theory, the DBM should have been able to find the building blocks for HIFF on the lower layer of its representation, and recombine them on the higher layers.

We now look at the results for solving at least 90% of runs (right three result columns of table 1). With the exception of the small onemax instances, and the concatenated 5-Traps problem, DBM-EDAs performance is inferior to BOA. In addition to the larger HIFF instance, DBM-EDA is unable to solve four of the NK landscape instances. In other words, while DBM-EDA was relatively competitive when the goal was to solve at least 50% of the instances to optimality, it seems to be less able to provide *reliable* convergence to the global optimum. A drastic example is the 150 bit instance of the simple onemax problem: DBM-EDA needs a population size of 200 to find the optimum in at least 50% of runs, but 6000 to find it in at least 90% of the runs. This behavior has also been observed in another EDA based on neural networks (DAE-EDA, see [18]⁷).

Second, we look at the CPU times required to solve the problem. For most instances, DBM-EDA is faster than BOA. Note that the direct comparison of CPU times is not entirely fair for BOA. In a more efficient programming language instead of a script-based language like Matlab/Octave, BOA’s speedup is significantly higher than the one of DAE-EDA. However, neither is Matlab/Octave the best programming language for DBM-EDA: Almost every recent implementation of neural networks is parallelized on graphics processing units (GPU), which, in turn, speeds up training and sampling these models considerably (see e.g. [25, 7, 26]). Parallelizing multivariate EDAs such as BOA is well possible, however the speedups are often single- or double-digit, even on GPUs (see e.g. [14, 13]). In contrast, parallelizing EDAs using neural networks can make proper use of modern GPU hardware and yield very high speedups: [21] report speedups of up to 200×, against optimized

⁷Note that this behavior is not so pronounced in later versions of DAE-EDA, which use the same model, but a different parametrization of the learning phase (publication [20] in preparation).

Problem	Algo	Average results					
		Population size such that optimum is found					
		in $\geq 50\%$ of runs			in $\geq 90\%$ of runs		
		PopSize	Unique Evals	Time (sec)	PopSize	Unique Evals	Time (sec)
ONEMAX50	BOA	125	2,119 \pm 125	685 \pm 101	125	2,119 \pm 125	685 \pm 101
	DBM	100	1,700\pm98	442\pm81	100	1,700\pm98	442\pm81
ONEMAX75	BOA	125	2,787 \pm 158	2,182 \pm 321	125	2,787 \pm 158	2,182 \pm 321
	DBM	100	2,328\pm86	565\pm88	100	2,328\pm86	565\pm88
ONEMAX100	BOA	250	6,259 \pm 153	8,967 \pm 1,016	250	6,259 \pm 153	8,967 \pm 1,016
	DBM	200	5,592\pm195	1,641\pm306	200	5,592\pm195	1,641\pm306
ONEMAX150	BOA	250	7,698 \pm 270	26,867 \pm 3,380	250	7,698\pm270	26,867\pm3,380
	DBM	200	6,095\pm503	3,555\pm283	6,000	212,464 \pm 6,846	73,024 \pm 8,317
4-Traps 40 bit	BOA	500	8,682\pm429	1,894 \pm 323	1,000	13,673\pm758	2,728 \pm 297
	DBM	2,000	34,561 \pm 6,168	2,034 \pm 880	3,000	47,231 \pm 5,712	2,201\pm312
4-Traps 60 bit	BOA	500	12,152\pm518	6,797 \pm 927	1,000	20,236\pm1,362	10,604 \pm 1,707
	DBM	4,000	89,495 \pm 4,723	5,481\pm416	5,000	104,967 \pm 12,482	6,793\pm686
4-Traps 80 bit	BOA	1,000	26,377\pm780	26,871 \pm 3,906	2,000	43,777\pm1,695	43,935 \pm 4,994
	DBM	6,000	153,278 \pm 6,149	13,271\pm1,295	6,000	153,278 \pm 6,149	13,271\pm1,295
5-Traps 25 bit	BOA	1,000	11,032 \pm 877	1,023 \pm 245	1,500	14,924 \pm 1,028	1,384 \pm 211
	DBM	1,000	10,368 \pm 2,125	444\pm108	1,500	13,291 \pm 2,471	566\pm108
5-Traps 50 bit	BOA	3,000	47,904 \pm 3,120	20,199 \pm 2,704	3,000	47,904 \pm 3,120	20,199 \pm 2,704
	DBM	3,000	51,367 \pm 19,948	3,168\pm2,087	4,000	49,886 \pm 11,933	3,060\pm617
5-Traps 75 bit	BOA	4,000	90,802\pm2,712	86,908 \pm 8,345	6,000	119,044 \pm 4,353	119,275 \pm 15,826
	DBM	5,000	99,990 \pm 6,169	8,538\pm1,131	6,000	101,107\pm19,431	9,183\pm1,407
5-Traps 100 bit	BOA	6,000	151,231\pm3,207	284,456 \pm 27,058	8,000	190,011 \pm 4,664	355,140 \pm 25,659
	DBM	8,000	169,700 \pm 12,290	26,802\pm6,659	8,000	169,700\pm12,290	26,802\pm6,659
NK $n = 30$, $k = 4$, $i = 1$	BOA	500	9,820 \pm 874	1,364 \pm 274	2,000	32,015\pm3,094	4,590 \pm 1,044
	DBM	300	5,976\pm941	453\pm88	5,000	69,081 \pm 9,269	2,120\pm576
NK $n = 30$, $k = 4$, $i = 2$	BOA	2,000	37,883 \pm 3,120	6,753 \pm 1,551	4,000	67,939\pm6,649	13,360 \pm 3,445
	DBM	2,000	30,124\pm3,941	1,508\pm310	9,000	126,982 \pm 13,840	3,529\pm680
NK $n = 34$, $k = 4$, $i = 1$	BOA	500	11,685 \pm 788	1,896 \pm 382	1,000	21,546\pm1,860	3,603 \pm 625
	DBM	400	8,823\pm1,227	567\pm105	4,000	69,125 \pm 6,476	2,807\pm631
NK $n = 34$, $k = 4$, $i = 2$	BOA	2,000	41,260\pm3,544	9,178 \pm 1,885	5,000	88,321\pm9,272	20,377\pm5,123
	DBM	4,000	65,752 \pm 6,832	2,758\pm525	-	-	-
NK $n = 30$, $k = 5$, $i = 1$	BOA	250	5,835 \pm 867	787 \pm 221	500	11,221\pm644	1,565 \pm 260
	DBM	200	4,284\pm851	611\pm290	1,000	17,617 \pm 1,601	1,045\pm178
NK $n = 30$, $k = 5$, $i = 2$	BOA	500	12,122\pm1,456	1,584 \pm 350	2,000	41,641\pm5,157	6,831\pm1,541
	DBM	2,000	33,963 \pm 3,124	1,773 \pm 247	-	-	-
NK $n = 34$, $k = 5$, $i = 1$	BOA	6,000	130,026 \pm 7,306	34,671 \pm 4,232	16,000	307,171\pm24,227	86,846\pm15,431
	DBM	6,000	106,392\pm10,966	4,153\pm793	-	-	-
NK $n = 34$, $k = 5$, $i = 2$	BOA	13,000	245,051 \pm 30,834	63,222 \pm 16,462	16,000	300,058\pm42,914	84,266\pm22,365
	DBM	10,000	192,125\pm18,928	6,752\pm1,254	-	-	-
HIFF64	BOA	500	11,991\pm731	7,480 \pm 1,059	500	11,991\pm731	7,480 \pm 1,059
	DBM	3,000	64,073 \pm 4,283	7,697 \pm 843	3,000	64,073 \pm 4,283	7,697 \pm 843
HIFF128	BOA	1,000	35,477\pm1,537	99,782\pm10,278	1,500	51,008\pm2,387	137,617\pm12,151
	DBM	-	-	-	-	-	-

TABLE 1. This table shows average results for fitness evaluations and CPU time for DBM-EDA and BOA for the test problems. For each instance and algorithm, we selected the minimal population size which leads to the optimal solution in at least 10 (left three result columns) or 18 (right three result columns) of 20 runs. Bold results are significantly smaller, according to a Wilcoxon signed-rank tests ($p < 0.01$, data is not normally distributed)

CPU code, for RBM-EDA, which uses a neural network model that is closely related to the DAE. Hence, it is reasonable to assume that an efficient GPU-based implementation of DBM-EDA will still be fast.

However, other neural network based EDAs such as DAE-EDA are computationally less expensive. Hence, they need considerably less time for solving the same benchmark instances to optimality ([18, 19, 20]).

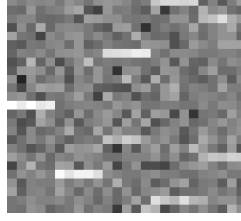


FIGURE 2. Visualization of the learned **weight matrix** \mathbf{W}^1 of DBM-EDA optimizing a concatenated 5-trap problem with 30 bits, in a later EDA generation.

4. DISCUSSION AND CONCLUSION

We introduced DBM-EDA, an Estimation of Distribution Algorithm which uses a Deep Boltzmann Machine as its probabilistic model. We evaluated the performance of DBM-EDA on multiple instances of standard benchmark problems for combinatorial optimization, and compared the results to the state-of-the-art multivariate Bayesian Optimization Algorithm.

DBM-EDA was able to solve most of the instances. However, its model quality was not competitive to BOA, with the exception of concatenated 5-Trap problems. Specifically, DBM-EDA was often unable to provide reliable convergence to the global optimum, or needed very large population sizes. Correspondingly, a high number of fitness evaluations were required. A reason for the insufficient model quality might be the mean-field approximation of the DBM’s training process. Mean-field approximation struggles if the probability distribution being approximated is multimodal. However, this is often the case in the early EDA generations: A sample might resemble configurations of different local optima, perturbed by noise. Surprisingly, DBM-EDA was unable to solve the larger HIFF instance at all. This is despite the fact that, as a hierarchical model with multiple layers, DBM-EDA should be particularly well-suited for a hierarchical optimization problem like HIFF.

DBM-EDA’s performance the concatenated trap problem with traps of size $k = 5$ was superior to BOA. Here, DBM-EDA was able to solve the larger instances (75 and 100 bit) with fewer fitness evaluations. Recall that the problem is particularly difficult, as it is composed of subproblems which are deceptive. We hypothesize that the reason for the good performance is the structure of the DBM: The hidden neurons are conditionally independent (see Equations 4 and 5). This matches the fitness function, which is additively decomposable, and separable. Figure 2 shows that neurons in the first hidden layer tend to model global optima to different additive parts of the fitness function. White pixels indicate large positive weight values, black pixels large negative weight values. Each row visualizes the connections between a single hidden neuron (of the first hidden layer) and the 30 problem variables. In the deceptive 5-trap problem, blocks of five adjacent variables have a strong contribution to the fitness, if *all five* variables are equal to one or equal to zero. Each block of five variables is independent of all other blocks. The figure shows that many hidden neurons strongly influence a single block of problem variables (bright/dark blocks of five adjacent pixels), and are indifferent to most other neurons (mid gray values). The learned representation of the model therefore

matches the problem structure. As the hidden neurons are activated stochastically, samples can comprise local optima of different additive terms of the fitness function, even if this specific combination has not been seen in the training data. A similar behavior has been observed for DAE-EDA, another EDA using neural networks [20].

In sum, while it is feasible to use a DBM as an EDA model, the effort for learning the multi-layered DBM model seems not to pay off for the optimization process in a noisy environment. There are multiple areas for future research. In the case where only 50% of the runs were required to find the global optimum, the results for DBM-EDA were quite encouraging. DBM-EDA could be a useful tool, if it could provide reliable convergence. The reasons to why this is currently not the case should be analyzed properly. Also, more work is necessary to understand why DBM-EDA was unable to apply the benefits of its hierarchical model to the HIFF problem.

REFERENCES

- [1] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski. A learning algorithm for boltzmann machines. *Cognitive science*, 9(1):147–169, 1985.
- [2] Y. Bengio. Learning deep architectures for AI. *Foundations and Trends Machine Learning*, 2(1):1–127, Jan. 2009.
- [3] K. Deb and D. E. Goldberg. Analyzing deception in trap functions. In D. L. Whitley, editor, *Foundations of Genetic Algorithms 2*, pages 93–108. Morgan Kaufmann, 1993.
- [4] G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- [5] G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- [6] S. A. Kauffman and E. D. Weinberger. The NK model of rugged fitness landscapes and its application to maturation of the immune response. *Journal of Theoretical Biology*, 141(2):211–245, 1989.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [8] P. Larrañaga and J. A. Lozano. *Estimation of Distribution Algorithms: A New Tool for Evolutionary Computation*, volume 2 of *Genetic Algorithms and Evolutionary Computation*. Springer US, Boston, MA, USA, 2002.
- [9] Y. LeCun, Y. Bengio, and G. E. Hinton. Deep learning. *Nature*, 521:436–444, 2015.
- [10] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, pages 609–616, New York, NY, USA, 2009. ACM.
- [11] B. L. Miller and D. E. Goldberg. Genetic algorithms, tournament selection, and the effects of noise. *Complex Systems*, 9:193–212, 1995.
- [12] H. Mühlenbein and G. Paaß. From recombination of genes to the estimation of distributions i. binary parameters. In H.-M. Voigt, W. Ebeling, I. Rechenberg, and H.-P. Schwefel, editors, *Parallel Problem Solving from Nature (PPSN IV)*, volume 1141 of *Lecture Notes in Computer Science*, pages 178–187. Springer, Berlin, Heidelberg, 1996.
- [13] A. Munawar, M. Wahib, M. Munetomo, and K. Akama. Theoretical and empirical analysis of a GPU-based parallel bayesian optimization algorithm. In *International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT 2009)*, pages 457–462. IEEE, 2009.
- [14] J. Očenášek and J. Schwarz. The parallel bayesian optimization algorithm. In P. Sinčák, J. Vaščák, V. K., and R. Mesiár, editors, *The State of the Art in Computational Intelligence*, volume 5 of *Advances in Soft Computing*, pages 61–67. Physica-Verlag, Heidelberg, 2000.
- [15] M. Pelikan. Bayesian optimization algorithm. In *Hierarchical Bayesian Optimization Algorithm*, volume 170 of *Studies in Fuzziness and Soft Computing*, pages 31–48. Springer, 2005.

- [16] M. Pelikan. Analysis of estimation of distribution algorithms and genetic algorithms on NK landscapes. In C. Ryan and M. Keijzer, editors, *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2008)*, volume 10, pages 1033–1040, New York, NY, USA, 2008. ACM.
- [17] M. Pelikan, D. E. Goldberg, and E. Cantu-Paz. BOA: the bayesian optimization algorithm. In W. Banzhaf, J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, and R. E. Smith, editors, *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 1999)*, volume 1, pages 525–532, San Francisco, CA, USA, 1999. Morgan Kaufmann.
- [18] M. Probst. Denoising autoencoders for fast combinatorial black box optimization. Technical Report arXiv:1503.01954, University of Mainz, 2015.
- [19] M. Probst. Denoising autoencoders for fast combinatorial black box optimization. In *Proceedings of the Companion Publication of the 2015 Genetic and Evolutionary Computation Conference (GECCO)*, GECCO Companion '15, pages 1459–1460, New York, NY, USA, 2015. ACM.
- [20] M. Probst and F. Rothlauf. Model building and sampling in estimation of distribution algorithms using denoising autoencoders. Technical report, University of Mainz, 2016.
- [21] M. Probst, F. Rothlauf, and J. Grah. An implicitly parallel EDA based on restricted boltzmann machines. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2014)*, pages 1055–1062, New York, NY, USA, 2014. ACM.
- [22] M. Probst, F. Rothlauf, and J. Grah. Scalability of using restricted boltzmann machines for combinatorial optimization. *to appear in: European Journal of Operational Research*, 2016.
- [23] R. Salakhutdinov and G. E. Hinton. Deep boltzmann machines. In D. van Dyk and M. Welling, editors, *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS 2009)*, volume 5, pages 448–455, 2009.
- [24] R. Salakhutdinov and G. E. Hinton. A better way to pretrain deep boltzmann machines. In *Advances in Neural Information Processing Systems*, pages 2447–2455, 2012.
- [25] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [26] I. Sutskever, O. Vinyals, and Q. V. V. Le. Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3104–3112. Curran Associates, Inc., 2014.
- [27] R. A. Watson, G. S. Hornby, and J. B. Pollack. Modeling building-block interdependency. In *Parallel Problem Solving from Nature - PPSN V*, pages 97–106. Springer, 1998.

JOHANNES GUTENBERG-UNIVERSITÄT MAINZ, DEPT. OF INFORMATION SYSTEMS AND BUSINESS
ADMINISTRATION, JAKOB-WELDER-WEG 9, 55128 MAINZ, GERMANY

E-mail address: {probst|rothlauf}@uni-mainz.de

URL: <http://wi.bwl.uni-mainz.de>